

Hugo Matos de Rangel Tavares

Testing Deterministic Regular Expressions



Departamento de Ciência de Computadores
Faculdade de Ciências da Universidade do Porto
Janeiro de 2020

Hugo Matos de Rangel Tavares

Testing Deterministic Regular Expressions

Curricular Project Report

Supervisor: Nelma Moreira and Rogério Reis

Departamento de Ciência de Computadores
Faculdade de Ciências da Universidade do Porto
Janeiro de 2020

Contents

1	Introduction	4
2	Elementary Definitions	5
2.1	Regular Expressions	5
2.2	DFA	6
2.3	NFA	6
2.4	Regular Expressions as Trees	7
2.4.1	LCA	8
2.4.2	LSA	8
3	Glushkov Automaton	9
4	Testing Determinism	12
4.1	Skeleton tree	13
4.2	BuildNext	15
4.3	CheckNode	16
4.4	isDeterministic	17
5	Conclusions	18
	Appendices	19
	References	27

Chapter 1

Introduction

Regular expressions are widely used for matching. Thus it is important to have efficient matching algorithms of a word against a regular expression. One way would be to convert the regular expression into a deterministic finite automaton (DFA) but this automaton can be exponentially larger than the expression. In alternative, nondeterministic finite automata (NFA) equivalent to a given expression can be at most quadratically larger.

However testing membership for a NFA is more costly. In this context, deterministic regular expressions are expressions that when simulated by Glushkov automaton the corresponding automaton is deterministic [1]. These expressions are most interesting and are in general used in applications such as XML schemas or network intrusion detection. One problem then is to know, given a general regular expression if it is deterministic. The conversion to Glushkov Automaton is quadratic in the size of the expression. Recently a linear time test was developed [3]. In this work we studied this method and implemented it using the `fado.dcc.fc.up.pt`, a library for manipulation of formal languages representations. This report is organized as follows.

First we describe some fundamental definitions such as regular expressions and both deterministic and non-deterministic finite automata, DFA and NFA respectively. Next we explain the construction method of Glushkov automaton and the sets that define it. Then, we present a method that checks whether a regular expression is deterministic or not that runs in $O(|e|)$, with e as the length of the tested regular expression. In the appendix we exhibit a Python 2.7 program with a complete tree structure for regular expressions, its construction and a few auxiliary functions to test determinism in a regular expression.

Chapter 2

Elementary Definitions

2.1 Regular Expressions

A regular expression, or regex, is a sequence of characters and rules that define a search pattern. A string matches a certain regex if and only if you can build the same string with the regex pattern.

Having a finite set of symbols (which we call alphabet denoted by Σ), a word is a sequence of symbols of Σ . For instance, $\Sigma = \{a,b\}$, then $aabb$ is a word. The ε word is the empty sequence. The set of every words is denoted by Σ^* .

A language L is a subset of Σ^* , $L \subseteq \Sigma^*$. A regular expression over Σ defines a language and the set of regular expressions R is:

- if $a \in R$, a is a regular expression and $L(a) = \{a\}$
- $\varepsilon \in R$, $L(\varepsilon) = \{\varepsilon\}$
- $e, e' \in R$, iff e is a regular expression and e' is a regular expression, then $e+e'$ and ee' are regular expressions in R .
- $L(e+e') = L(e) \cup L(e')$.
- $L(ee') = L(e) \cdot L(e') = \{xy \mid x \in e \text{ and } y \in e'\}$.
- if e is a regular expression, then e^* is a regular expression and $L(e^*) = (L(e))^*$.
- a regex can use other Boolean operators or the *option* operator such that $L(e?) = L(e+\varepsilon)$.

2.2 DFA

A deterministic finite automaton, DFA, is a finite state machine that tests whether a word belongs, or not, to the language described by the regular expression. If in the end of the word the current state is a final state (q is a final state iff $q \in F$) the word belongs to the language.

A DFA is described as a tuple with:

- Q , a finite set of states.
- Σ , an alphabet, i.e. a finite set of symbols.
- $\delta: Q \times A \rightarrow Q$, a transition function, the term deterministic comes from the fact that each transition takes a current state and a symbol has arguments and returns a next state, always and only one for each tuple (state,symbol).
- $q_0 \in Q$, a start state.
- $F \subseteq Q$, a set of final states.

Therefore, we represent the automaton A as a quintuple: $A = (Q, \Sigma, \delta, q_0, F)$.

2.3 NFA

A non-deterministic finite automata (NFA) is an automaton that also describes the acceptance of a language but it has the possibility *to be in several states at once*. This ability is often expressed as an ability to estimate something about its input, for instance, when we use an automaton to search for a certain word in a long text string it can be essential to guess the current position we are in the string that we search using a chain of states to check if the string exists.

Similarly to a DFA, an NFA has a finite set of states, a finite set of symbols *alphabet*, a starting state and a set of accepting states (*final states*). They differ on the transition function where the NFA can have $\delta(X, 'a')=Y$ and $\delta(X, 'a')=Z$ with $Y \neq Z$, but DFA can not. For the NFA, the transition function takes a state and an input symbol as arguments and can return zero, one or more states, not like the DFA which can only return one, i.e. $\delta \subseteq Q \times A \times Q$.

2.4 Regular Expressions as Trees

In order to test determinism of a regular expression is possible to represent the regex as a tree structure where the nodes are the operators and the leaves are the symbols of the alphabet. To build this tree we have to define:

- $node(regex, path, tree)$, this constructor creates a node in the tree where $regex$ is the regular expression denoted in that node, the $path$ is a string that represents the path from the root ($path = '1'$) to that node and the $tree$ that owns the node.
- *NodeTable* : This table keeps records of a path to each node for an easier access to a node. In the Python implementation, *NodeTable* represented as a dictionary having paths as keys and nodes as values so that we can get a node from the path. For example, if we are in a concatenation node if we want to get the left child (*Lchild*) we just need to search in our *NodeTable* for the path of the current node and add a '1' to the path, or a '2' if we want the right child (*Rchild*).

In the tree construction, a node has a left child if is labeled as star, a left and a right child if it is a concatenation or disjunction or no childs if it is a symbol, in this case we add the symbol to our tree's *alphabet* as shown in *Algorithm 1*. The procedure *Construct* receives a node that expresses an expression and first will build a node for it, then it looks to the expression type and if it is a symbol, then adds it to the language's alphabet, if it is a star expression then he builds it only child (expression without kleene star), otherwise constructs the left child then the right child of the concatenation or the disjunction.

Algorithm 1 Computing a regex tree.

```

procedure CONSTRUCT(regex,path, $t_e$ )
   $n \leftarrow Node(regex, path, t_e)$ 
   $NodeTable[path] \leftarrow n$ 
  if Atom( $n$ ) then
     $alphabet(t_e) \leftarrow alphabet(t_e) \cup symbol(n)$ 
  else
    if lab( $n$ ) is * then
       $Lchild(n) \leftarrow Construct(regex.arg, path +' 1', t_e)$ 
    else
       $Lchild(n) \leftarrow Construct(regex.arg1, path +' 1', t_e)$ 
       $Rchild(n) \leftarrow Construct(regex.arg2, path +' 2', t_e)$ 
  return  $n$ 

```

2.4.1 LCA

The *LCA* or lowest common ancestor, is a function that takes two nodes from our structure as arguments and returns the lowest common node in the tree. This can be seen has the smaller sub-expression that contains both sub-expressions (nodes) in the arguments. This function can be implemented in constant time regardless of the length of the expression. [3]

2.4.2 LSA

The *LSA*, lowest star ancestor, is a function that takes one node from a regex tree as argument and returns the first star-labeled ancestor that the algorithm finds by going through the tree from the argument to the root, if the *LSA* get to root node and this is not a star-labeled node then the argument has no *LSA*. [3] With a tree structure organized with sub-expressions in every node where the root node contains the entire regular expression, we can find the *LSA* of a node in $O(|e|)$ where $|e|$ denotes the length of the expression e .

Chapter 3

Glushkov Automaton

Given the fact that a regular expression defines a pattern of words from a regular language, we can define a set of words that belong to a certain language by neatly browsing the expression that establishes the language. For instance, we can obtain the word $aabaa$ from the expression $a*ba*$ by using twice the first a , then using b and then, again, twice the second a , resulting on a^2ba^2 . It is important to understand that a word can be described by more than one expression and also the other way around. Taking by example the word above, we can also say that $aabaa$ can be accepted by $((aa)*b)*$. It is known that, the order of the symbols are crucial in word recognition. From this fact, Victor Glushkov, idealized a construction of a NFA with no ε -transitions. [2]

Glushkov's automata are based on distinguish all symbols from the expression by giving them an index regardless their label, then take them as states and establishing transitions according the follow-up among the symbols.

Below we present the procedure that Glushkov's method implements to mark the regular expression.

$$\begin{aligned}mark(\varepsilon, i) &= \varepsilon \\mark(a, i) &= a_i \\mark(a + e, i) &= a_i + mark(e, i + 1) \\mark(a.e, i) &= a_i.mark(e, i + 1) \\mark(e*, i) &= mark(e, i)\end{aligned}$$

The marked symbols are called positions, and the set of positions is represented by $Pos(e)$.

Using one of the above examples, $a * ba*$, we can demonstrate the process by indexing all the symbols, getting $(a,1)^*(b,2)(a,3)^*$, and observe that $(a,1)$ is always

followed by either (a,1) or (b,2), (b,2) is followed by (a,3) or \emptyset and (a,3) by (a,3) or \emptyset . When a state can be "followed" by \emptyset it means that this is a final state. If the word ends there it belongs to the language described by the expression.

Indexing and identifying the symbols is not enough, we need a transition function to determine our automata's edges. So, Glushkov, added three sets to his construction: *First*, *Last* and *Follow*. The *First*(e) set contains all symbols that can show up in the beginning of the word. The *Last*(e) set contains all symbols that can appear in the end. The *Follow*(e) set, is a relation which we represent as a *dictionary* that attributes to every indexed symbol a set of symbols that can follow it. Next, we present how Glushkov's method generates both *First*, *Last* and *Follow* sets.

$$\begin{aligned} First(\varepsilon) &= First(\emptyset) = \emptyset \\ First(\sigma_i) &= \sigma_i \\ First(e_1 + e_2) &= First(e_1) \cup First(e_2) \\ First(e_1 e_2) &= First(e_1) \cup \varepsilon(e_1) First(e_2) \\ First(e^*) &= First(e) \end{aligned}$$

$$\begin{aligned} Last(\varepsilon) &= Last(\emptyset) = \emptyset \\ Last(a_i) &= a_i \\ Last(e_1 + e_2) &= Last(e_1) \cup Last(e_2) \\ Last(e_1 e_2) &= Last(e_1) \cup \varepsilon(e_2) Last(e_1) \\ Last(e^*) &= Last(e) \end{aligned}$$

$$\begin{aligned} Follow(\varepsilon, a_i) &= Follow(a_i, a_j) = \emptyset \\ Follow(e_1 + e_2, a_i) &= \begin{cases} Follow(e_1, a_i) & , a_i \in Pos(e_1) \\ Follow(e_2, a_i) & , a_i \in Pos(e_2) \end{cases} \\ Follow(e_1 \cdot e_2, a_i) &= \begin{cases} Follow(e_1, a_i) & , a_i \in Pos(e_1) \setminus Last(e_1) \\ Follow(e_1, a_i) \cup First(e_2) & , a_i \in Last(e_1) \\ Follow(e_2, a_i) & , a_i \in Pos(e_2) \end{cases} \\ Follow(e_1^*, a_i) &= \begin{cases} Follow(e_1, a_i) & , a_i \in Pos(e_1) \setminus Last(e_1) \\ Follow(e_1, a_i) \cup First(e_1, a_i) & , a_i \in Pos(e_2) \end{cases} \end{aligned}$$

Note that $\varepsilon(e_1)$ implies that e_1 accepts ε and in this case we say that e_1 is nullable.

With these sets we describe Glushkov's automaton as $AG = (Pos(e) \cup \{0\}, \Sigma, \delta, Last(e) \cup \varepsilon(0))$, where $\delta(0) = First(e)$ and $\delta(a_i) = Follow(e, a_i)$. The following Lemmas allow the efficient computation of the *Follow* set [3]. We denote by *tree*, the tree representation of the regular expression e .

Lemma 1 *Let $p, q \in Pos(tree)$, $n = LCA(p, q)$ and $s = LSA(p, q)$. Then $q \in Follow(p)$ iff:*

- (1) $lab(n) = concat$, $q \in First(Rchild(n))$, $p \in Last(Lchild(n))$.
- (2) $q \in First(s)$, $p \in Last(s)$.

We say that $q \in Follow^\bullet(p)$ iff (1) is satisfied, and $q \in Follow^*(p)$ iff (2) is satisfied. Note that is possible for some nodes p and q to satisfy both (1) and (2).

Now, we define the Boolean properties *SupFirst* and *SupLast* for every node n , where $n' = parent(n)$:

- *SupFirst*(n) iff $label(n') = concat$, $n = Rchild(n')$ and $Lchild(n')$ is not nullable.
- *SupLast*(n) iff $label(n') = concat$, $n = Lchild(n')$ and $Rchild(n')$ is not nullable.

With the *Last* set defined is possible to determine the set *FollowAfter* which indicates possible symbols to appear next to the sub-expression (*node*) we want. To do that we also need a function to determine if a node a is ancestor of a node b , so, we implemented *Reflexive*(a, b) that takes two arguments and returns true if b is ancestor of a and false otherwise. (see *Algorithm 2*)

Algorithm 2 Computing *Follow* of a sub-expression.

```

procedure FOLLOWAFTER(tree,node)
   $S \leftarrow \emptyset$ 
  for  $p$  in  $Last(node)$  do
    for  $q$  in  $Pos(tree)$  do
      if  $Follow(p, q)$  and  $\neg Reflexive(q, node)$  then
         $S \cup = q$ 
  return  $S$ 

```

Testing Determinism

It was shown before [1] that a regular expression is deterministic iff its Glushkov automaton is deterministic, i.e. a DFA.

In this chapter we describe a determinism test for regular expressions that avoids the construction of Glushkov automaton and can run in linear time.

In *Figure 4.1* we represent an example of a regex and the corresponding tree annotated with the several properties defined in the previous chapter.

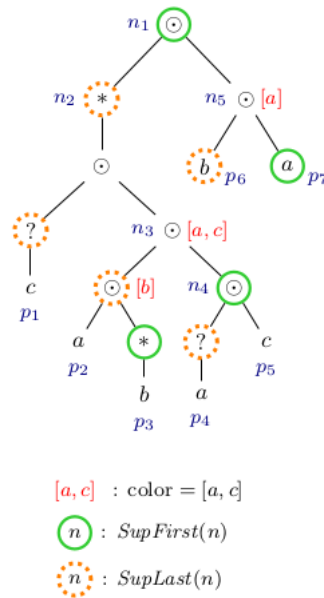


Figure 4.1: Expression $e = (c?((ab^*)(a?c)))^*(ba)$ represented in a tree. [3]

Let $q, q' \in Nodes(tree)$, where $q \neq q'$ and $n = qSupFirst(q) = qSupFirst(q')$.

With *First* and *Last* sets defined for every node in the tree, for $p \in \text{Last}(\text{Lchild}(n'))$ where n' is the parent of n we can say that $\text{parent}(n) = \text{LCA}(p, q) = \text{LCA}(p, q')$.

By Lemma 1, $q, q' \in \text{Follow}(p)$ and by the definition of determinism $\text{lab}(q) \neq \text{lab}(q')$.

Proposition 1 $\forall q \neq q' \in \text{Nodes}(\text{tree}), p\text{SupFirst}(q) = p\text{SupFirst}(q') \text{ implies } \text{lab}(q) \neq \text{lab}(q')$.

Knowing that a deterministic expression requires a deterministic automaton, if we test every node of the tree, the *Follow* set cannot have two different atoms with the same color.

Proposition 2 $\forall a \in \Sigma \text{ and } n \in \text{Nodes}(t_a), \text{ where } t_a \text{ is defined below, } \text{Next}(n, a) \text{ contains at most one element.}$

The function $\text{Next}(n, a)$ returns a set of atoms that belong to the set $\text{Follow}(n)$ labeled a .

In order to analyse the regex tree in linear time we establish a color set to each node of our tree meaning that from a position p colored a , the next letter from the word can be a symbol a . So we assign:

- $\forall p \in \text{Pos}(\text{tree})$ and $p' = p\text{SupFirst}(p)$, include $\text{lab}(p)$ in $\text{colors}(\text{parent}(p'))$.
- p is a *witness* for color a in p' .

Note that a node may have more than one color assigned to it, but taking into account *Proposition 1*, each node can have at most one witness per color to be deterministic. This implies that a node n cannot have as witness p_2 and p_4 for color a , this would mean that during the test when we analyse n if we read an a , the automaton would jump to states p_2 and p_4 defining the expression as non-deterministic.

4.1 Skeleton tree

The *skeleton tree* consist on all positions in t_e labeled a , their $p\text{SupLast}$ and LSA nodes, and every LCA of two nodes of class a , denoted in t_a .

In our implementation, we defined t_a as a set of pointers to nodes from t_e that also composes the *skeleton tree*. This tree, however, has a peculiarity, not like t_e the

skeleton tree can have a node labeled *concat* with only one son, being the second *None*.

In *Figure 4.2* we can observe that the *a-skeleton* from t_e generates an unambiguous tree that shows the precedence among symbols labeled *a*.

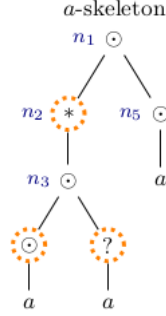


Figure 4.2: Skeleton for color *a* in the tree of *Figure 1*. [3]

Algorithm 3 Computing the *skeleton tree* for color *a*.

```

procedure buildta(tree, a)
  atoms(ta)  $\leftarrow \emptyset$ 
  for n in atoms(tree) do
    if lab(n) = a then
      atoms(ta)  $\cup = n$ 
  ta  $\leftarrow atoms(t_a)$ 
  for x in atoms(ta) do
    for y in atoms(ta) do
      if x = y then
        ta  $\cup = parent(x, tree)$ 
      else
        ta  $\cup = LCA(x, y)$ 
  ta  $\cup = find\_stars(tree, a)$ 
  return ta

```

The function *find_stars*(*tree*, *a*) returns a set with all *star*-labeled nodes that are ancestors of any position labeled *a* from our tree.

4.2 BuildNext

In order to test the determinism of our regex, we implemented a function *BuildNext* which tests determinism for all nodes in t_e . This function will check every node, one by one, and test if there is at least one node where $Next(n, a)$ has more than one element, or two if $lab(n) = +$. If it occurs in at least one node, then the regex is non-deterministic.

Algorithm 4 [3] Computing the set $Next(n, a)$.

```

procedure BUILDNEXT( $a, n, Y$ )
  if  $SupLast(n)$  then
     $Y \leftarrow \emptyset$ 
  if  $n$  is the left child in  $t_a$  of a concat node and
     $n$  has a right sibling  $n'$  in  $t_a$  and
     $(\neg SupLast(node) \text{ or } parent(n, t_a) = parent(n, t_e))$  then
       $Y \leftarrow Y \cup \{FirstPos(n', a)\}$ 
   $Next(n, a) \leftarrow \{p \in Y \mid n \prec p\}$ 
  if  $lab(n) = *$  then
     $Y \leftarrow Y \cup \{FirstPos(n, a)\}$ 
  if  $|Y| > 2$  then
    return false
   $l \leftarrow Lchild(n, t_a)$ 
  if  $l = \emptyset$  then
    return true
  else
     $B \leftarrow BuildNext(a, l, Y)$ 
   $r \leftarrow Rchild(n, t_a)$ 
  if  $r = \emptyset$  then
    return B
  else
    return  $B \wedge BuildNext(a, r, Y)$ 

```

Executing $BuildNext(n, a, \emptyset)$ the algorithm will add to our tree structure a new dictionary that will link a node n to the *positions* that can come next in the word. This dictionary will be represented as $Next$ and will be denoted as $Next(n)$ to represent $Follow(n)$, or $Next(n, a)$ a subset of $Follow(n)$ that only keeps the positions labeled a . With this we can also realize that $Next(n, a)$ can only have one element, or two if n is labeled $+$ (disjunction).

4.3 CheckNode

Here we define a function that verifies if a sub-expression (represented as a tree node) is deterministic. This function will be later tested for every node in our tree. All nodes (sub-expressions) have to be deterministic so that the expression e is deterministic, in other words $\forall n \in \text{nodes}(\text{tree}) \wedge \forall a \in \Sigma \text{checknode}(n, a)$ returns *true*. This function is defined in *Algorithm 5*.

Algorithm 5 Verify determinism in node *node*.

```

procedure CHECKNODE(node, sym)
  te  $\leftarrow$  tree(node)
  ta  $\leftarrow$  build $_t$ (te, sym)
  fp  $\leftarrow$  FirstPos(ta, n)
  s  $\leftarrow$  LSA(node)
  r  $\leftarrow$  Rchild(node)
  if  $\varepsilon(r.exp)$  then
    if has_next(te, n, sym) then return false
    if s is None then return true
    sl  $\leftarrow$  pSupLast(node)
    if FirstPos(ta, s) is fp and Reflexive(s, sl) then return false
  return true

```

With this function we can say that if we apply *checkNode* for n and his right child is non-nullable then we can say that n is deterministic because there is no possible next symbol in that sub-expression.

If the right child of n is nullable then we need to verify if the sub-expression in n can read the symbol a next, and also if in the path from $LSA(n)$ to $FirstPos(n, a)$ there is nothing non-nullable in the *left* childs and consequently that $FirstPos(n, a)$ follows the same position p that $Witness(n, a)$ follows.

Algorithm 6 Testing determinism for every node in t_e

```

procedure CHECKDETERMINISM( $tree, n$ )
  for  $x$  in  $colors(n)$  do
    if  $\neg checkNode(n, x)$  then return false
  if  $Lchild(n, tree)$  is  $None$  then return true
  else
     $l \leftarrow checkDeterminism(tree, Lchild(n, tree))$ 
  if  $Rchild(n, tree)$  is  $None$  then return  $l$ 
  else
     $r \leftarrow checkDeterminism(tree, Rchild(n, tree))$ 
  return  $l \wedge r$ 

```

With $checkDeterminism(tree, n)$ (defined in *Algorithm 6*) we first call it with the $root(tree)$ node so that the algorithm tests the function for every node from top to bottom recursively with a depth first search pattern.

This function first calls the function $checkNode$ on n for every color assigned to n , and if any call to the function $checkNode$ returns *false* than e is non-deterministic, otherwise $checkDeterminism$ will be tested for the childs of n , first to the left like DFS search.

4.4 isDeterministic

In order to verify the determinism of a regex e we test *Proposition 1* and *Proposition 2*, if any returns *False* then e is not deterministic. If both return *True*, then we execute $CheckNode(n, a)$ for every $a \in \Sigma$ and every node n coloured a . If any call to the function $CheckNode$ returns *False* then e is not deterministic.

This test is defined in *Algorithm 7*.

Algorithm 7 Testing determinism for a regex tree.

```

1: procedure ISDETERMINISTIC( $tree$ )
2:    $r \leftarrow root(tree)$ 
3:   for  $x$  in  $\Sigma(tree)$  do
4:      $buildNext(x, r, \emptyset)$ 
5:   if  $\neg P1(tree)$  or  $\neg P2(tree)$  then return false
   return  $checkDeterminism(tree, r)$ 

```

Chapter 5

Conclusions

In this project we have revised the Glushkov's automaton construction. Then, we presented a linear time algorithm to test if a regular expression is deterministic using properties acquired from Glushkov's construction, e.g. *First*, *Last* and *Follow* sets.

In the appendix we present an implementation of the method studied along with the implementation of a tree structure used in the test for determinism. It could not go without saying that FAdo System made possible the usage of a object oriented structure to assist during the construction of the tree and during the testing by facilitating the access to attributes of each regular expression.

We also performed several experiments of regular expressions generated uniform randomly. In order for this experiences be statistically significant further studies are needed. We also need to implement some of the key functions more carefully in order that the current implementation of the test for deterministic is indeed performed in linear time.

To complete, we note that the linear algorithm for testing determinism extends to regular expressions with numeric occurrence indicators, and that determinism can be decided in linear time.

Appendices

Here we present the code implemented in Python 2.7 within the FAdo system.

```

1 from os.path import commonprefix
2 from FAdo.reex import *
3 from FAdo.cfg import *
4 from FAdo.fa import *
5
6 # Tree Structure
7 class eTree:
8     def __init__(self, regex):
9         ReExp = concat(atom("#"), concat(regex.marked(), atom("$")))
10        self.NT = dict()
11        self.atoms = set()
12        self.next = dict()
13        self.sub_trees = dict()
14        self.alpha = set()
15        self.full_tree = construct(ReExp, "1", self)
16        self.root = self.NT.get("121")
17        self.Color()
18
19    def followList(self):
20        l = dict()
21        for x in self.atoms:
22            f_list = set()
23            for w in self.atoms:
24                if w.type != epsilon:
25                    if w.follow(x):
26                        f_list.add(w)
27            l[x] = f_list
28        return l
29
30    def build_ta(self, a):
31        if self.sub_trees.get(a) is not None:
32            return self.sub_trees.get(a)
33        ta_atoms = set()
34        for n in self.atoms:
35            if n.exp.val[0] == a:
36                ta_atoms.add(n)
37        ta = set()
38        for x in ta_atoms:
39            for y in ta_atoms:
40                if x is y:
41                    ta.add(x.parent)
42                else:
43                    ta.add(self.NT.get(commonprefix([x.path, y.path])))
44        self.sub_trees[a] = ta_atoms.union(ta.union(self.find_stars(a)))
45        return self.sub_trees.get(a)
46
47    def find_stars(self, a):

```

```

48     ta = set()
49     for n in self.atoms:
50         if n.exp.val[0] == a:
51             n = n.parent
52             while n != self.root:
53                 if n.type is star:
54                     ta.add(n)
55                     n = n.parent
56     return ta
57
58 def has_Next(self, n, a):
59     s = self.next.get(n)
60     if s is None:
61         return False
62     for x in s:
63         if x.exp.val[0] == a:
64             return True
65     return False
66
67 def Color(self):
68     for x in self.atoms:
69         x.pSupFirst().parent.colors.add(x.exp.val[0])
70     return
71
72 def isDeterministic(self):
73     for x in self.alpha:
74         buildNext(x, self.root, set())
75     if not self.cond_P1() or not self.cond_P2():
76         return False
77     return self.check_determinism(self.root)
78
79 def cond_P1(self):
80     for a in self.atoms:
81         for b in self.atoms:
82             if a is b:
83                 continue
84             if a.pSupFirst() == b.pSupFirst():
85                 if a.exp.val[1] == b.exp.val[1]:
86                     print "P1"
87                     return False
88     return True
89
90 def cond_P2(self):
91     for a in self.alpha:
92         ta = self.build_ta(a)
93         for n in ta:
94             k = n.FollowAfter
95             count = 0
96             for x in k:

```

```

97         if x.exp.val[0] == a:
98             count = count + 1
99         if count > 1:
100             print "P2"
101             return False
102     return True
103
104     def check_determinism(self, n):
105         for x in n.colors:
106             if not checkNode(n, x):
107                 return False
108         if n.left is not None:
109             l = self.check_determinism(n.left)
110         else:
111             return True
112         if n.right is not None:
113             return (l and self.check_determinism(n.right))
114         return l
115
116 # Searches for non-determinism in each node
117 def checkNode(n, a):
118     te = n.tree
119     ta = te.build_ta(a)
120     f = n.FirstPos(ta)
121     s = n.LSA()
122     if n.right.exp.ewp():
123         if te.has_Next(n, a):
124             return False
125         if s is None:
126             return True
127         if s.FirstPos(ta) is f and s.reflexive(n.pSupLas()):
128             return False
129     return True
130
131 # Constructor for eTree structure. Construct nodes.
132 def construct(reg_exp, path, t):
133     node = Node(reg_exp, path, t)
134     t.NT[node.path] = node
135     if node.type is star:
136         node.left = construct(node.exp.arg, path+'1', t)
137     elif node.type is concat or node.type is disj:
138         node.left = construct(node.exp.arg1, path+'1', t)
139         node.right = construct(node.exp.arg2, path+'2', t)
140     elif node.type is position:
141         sym = node.exp.val[0]
142         t.alpha.add(sym)
143     return node
144
145 class Node:

```

```

146     def __init__(self,expr,string,t):
147         self.type = type(expr)
148         self.path = string
149         self.exp = expr
150         self.tree = t
151         self.right = None
152         self.left = None
153         self.first = set()
154         self.last = set()
155         self.lsa = None
156         self.supfirst = None
157         self.suplast = None
158         self.FollowAfter = set()
159         self.colors = set()
160         if self.type is position:
161             t.atoms.add(self)
162         if string == '1':
163             self.parent = None
164         else:
165             self.parent = t.NT.get(string[:-1])
166
167     # check if n is ancestor of self
168     def reflexive(self,n):
169         if n is None:
170             return False
171         if commonprefix([self.path,n.path])==n.path: #LCA(self,n)
172             return True
173         return False
174
175     # boolean that states if an atom is the first from the right child
176     # of a concatenation
177     def SupFirst(self):
178         father = self.parent
179         if father is not None and father.type is concat:
180             if father.right is self and not father.left.exp.ewp():
181                 return True
182             return False
183
184     # boolean that states if an atom is the last from the left child of
185     # a concatenation
186     def SupLast(self):
187         father = self.parent
188         if father is not None and father.type is concat:
189             if father.left is self and not father.right.exp.ewp():
190                 return True
191             return False
192
193     # pointer to the first(right child) of a concatenation
194     def pSupFirst(self):

```

```

193         if self.supfirst is not None:
194             return self.supfirst
195         self.supfirst = None
196         if self.SupFirst() or self.path=="1":
197             self.supfirst = self
198         else:
199             self.supfirst = self.parent.pSupFirst()
200         return self.supfirst
201
202     # pointer to the last(left child) of a concatenation
203     def pSupLast(self):
204         if self.suplast is not None:
205             return self.suplast
206         if self.SupLast() or self.path=="1":
207             self.suplast = self
208         else:
209             self.suplast = self.parent.pSupLast()
210         return self.suplast
211
212     #commonprefix is the Python prelude function used for LCA
213     def follow(self,p):
214         node1 = self.tree.NT.get(commonprefix([p.path,self.path]))
215         if node1 is None:
216             return False
217         if node1.type is concat:
218             if self.follow_concat(p,node1):
219                 return True
220         node = node1.LSA()
221         if node is not None:
222             if self.follow_star(p,node):
223                 return True
224         return False
225
226     # Follow by Concatenation
227     def follow_concat(self,p,lca):
228         if self in lca.right.First() and p in lca.left.Last():
229             return True
230         return False
231
232     # Follow by Star
233     def follow_star(self,p,lsa):
234         if self in lsa.First() and p in lsa.Last():
235             return True
236         return False
237
238     def First(self):
239         if self.first != set():
240             return self.first
241         for node in self.tree.atoms:

```



```

242         if isFirst(self,node):
243             self.first.add(node)
244         return self.first
245
246     def Last(self):
247         if self.last != set():
248             return self.last
249         for node in self.tree.atoms:
250             if isLast(self,node):
251                 self.last.add(node)
252         return self.last
253
254     def LSA(self):
255         if self.lsa is not None:
256             return self.lsa
257         if self.type is star:
258             self.lsa = self
259         elif len(self.path) == 1:
260             self.lsa = None
261         else:
262             self.lsa = self.parent.LSA()
263         return self.lsa
264
265     def FirstPos(self,ta):
266         k = set()
267         for x in ta:
268             if x in self.First():
269                 k.add(x)
270         return k
271
272     def follow_after(self):
273         s = set()
274         for p in self.Last():
275             for q in self.tree.atoms:
276                 if q.follow(p) and not q.reflexive(self):
277                     s.add(q)
278         self.FollowAfter = s
279         return s
280
281     def isFirst(n,p):
282         k = p.pSupFirst()
283         if k is not None:
284             if p.reflexive(n) and n.reflexive(k):
285                 return True
286         return False
287
288     def isLast(n,p):
289         k = p.pSupLast()
290         if k is not None:

```

```

291         if p.reflexive(n) and n.reflexive(k):
292             return True
293     return False
294
295 def buildNext(a,n,y):
296     if n.SupLast():
297         y = set()
298     te = n.tree
299     ta = n.tree.build_ta(a)
300     if ta is None:
301         return False
302     dad = get_parent(n,ta)
303     if dad is not None:
304         r = get_right(dad,ta)
305         if dad.type is concat and get_left(dad,ta) is n and r is not
None:
306             if not n.SupLast() or get_parent(n,ta) is n.parent:
307                 y = y.union(r.FirstPos(ta))
308             n.tree.next[n] = set()
309             for x in y:
310                 if not x.reflexive(n):
311                     n.tree.next[n].add(x)
312             if n.type is star:
313                 y = y.union(n.FirstPos(ta))
314             if len(y)>2:
315                 return False
316             l = get_left(n,ta)
317             if l is None or l.exp.ewp():
318                 return True
319             else:
320                 b = buildNext(a,l,y)
321                 r = get_right(n,ta)
322                 if r is None or r.exp.ewp():
323                     return b
324                 return (b and (buildNext(a,r,y)))
325
326 def get_parent(n,tree):
327     t = n.tree
328     while n is not t.root:
329         parent = n.parent
330         if parent in tree:
331             return parent
332         n = n.parent
333     return None
334
335 def get_left(n,tree):
336     if n.type is position:
337         return None
338     path = n.path+'1'

```

```

339     possible = set()
340     for x in tree:
341         if len(x.path) >= len(path) and commonprefix([x.path, path]) == path:
342             possible.add(x.path)
343     if len(possible) == 0:
344         return None
345     return n.tree.NT.get(min(possible))
346
347 def get_right(n, tree):
348     if n.type is position:
349         return None
350     path = n.path + '2'
351     possible = set()
352     for x in tree:
353         if len(x.path) >= len(path) and commonprefix([x.path, path]) == path:
354             possible.add(x.path)
355     if len(possible) == 0:
356         return None
357     return n.tree.NT.get(min(possible))

```

References

- [1] Anne Brüggemann-Klein. Regular expressions into finite automata. *Theor. Comput. Sci.*, 120(2):197–213, 1993.
- [2] Hugo Gouveia. Obtenção de autómatos finitos não determinísticos pequenos. *UC Projeto*, 2009.
- [3] Benoît Groz and Sebastian Maneth. Efficient testing and matching of deterministic regular expressions. *J. Comput. Syst. Sci.*, 89:372–399, 2017.